

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平11-249679

(43) 公開日 平成11年(1999) 9月17日

(51) Int.Cl.⁶

G 1 0 L 3/00
5/04

識別記号

F I

G 1 0 L 3/00
5/04

H
F

審査請求 未請求 請求項の数 5 O L (全 5 頁)

(21) 出願番号 特願平10-52361

(22) 出願日 平成10年(1998) 3月4日

(71) 出願人 000006747

株式会社リコー

東京都大田区中馬込1丁目3番6号

(72) 発明者 酒寄 哲也

東京都大田区中馬込1丁目3番6号 株式
会社リコー内

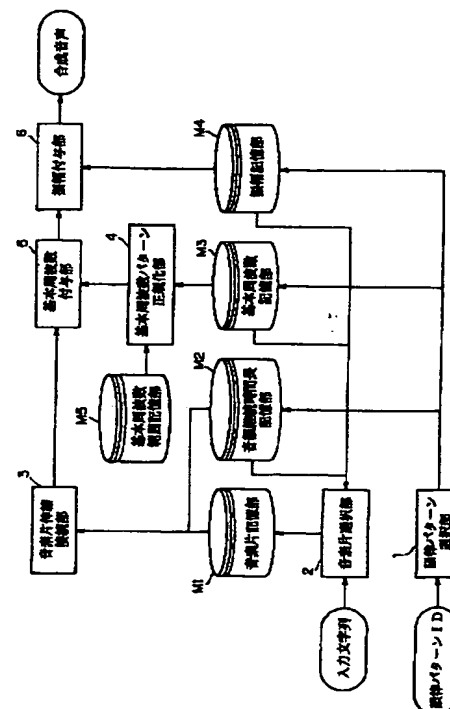
(74) 代理人 弁理士 高野 明近

(54) 【発明の名称】 音声合成装置

(57) 【要約】

【課題】 入力された文字列から定型的なフレーズを合成音声で読み上げる際の合成音声を、より自然性の音声を近いものにする。

【解決手段】 人間が発声したフレーズから抽出した音韻継続時間長系列、基本周波数系列、振幅あるいはパワーの系列をそれぞれ記憶した、音韻継続時間記憶部、基本周波数記憶部、振幅記憶部を備え、かつ、音素あるいは音素連鎖を音韻単位とし音韻情報を音素片としてを記憶する音素片記憶部を備えた音声合成装置により、入力文字列にしたがって音素片記憶部から読み出した音素片系列を、韻律パターン記憶部から読み出した音韻継続時間長系列にしたがって伸縮して接続し、韻律パターン記憶部から読み出した基本周波数系列にしたがって基本周波数付与を行い、その際、基本周波数範囲内に収まるように正規化を施し、更に、韻律パターン記憶部から読み出した振幅又はパワー系列にしたがって振幅付与を行って音声を合成する。



【特許請求の範囲】

【請求項1】 人間が発声したフレーズから抽出した音韻継続時間長系列、基本周波数系列、振幅あるいはパワーの系列をそれぞれ記憶した、音韻継続時間記憶部、基本周波数記憶部、振幅記憶部、及び、音素あるいは音素連鎖を音韻単位とし音韻情報を音素片として記憶する音素片記憶部を具備し、入力文字列にしたがって音素片記憶部から読み出した音素片系列を、韻律パターン記憶部から読み出した音韻継続時間長系列にしたがって伸縮して接続し、韻律パターン記憶部から読み出した基本周波数系列にしたがって基本周波数付与を行い、韻律パターン記憶部から読み出した振幅又はパワー系列にしたがって振幅付与を行って音声を作成することを特徴とする音声合成装置。

【請求項2】 請求項1に記載された音声合成装置において、音素片記憶部に記憶した音素片セットによって無理なく合成できる基本周波数の範囲を記憶する基本周波数範囲記憶部を具備し、基本周波数記憶部から読み出した基本周波数系列に対して、基本周波数範囲記憶部から読み出した基本周波数範囲内に収まるように正規化を施して音声を作成することを特徴とする音声合成装置。

【請求項3】 請求項1に記載された音声合成装置において、前記音素片記憶部は一つの音韻単位に対して適応すべき基本周波数範囲毎に複数の音素片を記憶しており、前記基本周波数記憶部から読み出した基本周波数に対応した音素片を選択的に用いて音声合成を行うことを特徴とする音声合成装置。

【請求項4】 請求項1に記載された音声合成装置において、前記音素片記憶部は一つの音韻単位に対して適応すべき振幅範囲毎に複数の音素片を記憶しており、基本周波数記憶部から読み出した振幅あるいはパワーに対応した音素片を選択的に用いて音声合成を行うことを特徴とする音声合成装置。

【請求項5】 請求項1に記載された音声合成装置において、前記音素片記憶部は一つの音韻単位に対して適応すべき音韻継続時間長範囲毎に複数の音素片を記憶しており、音韻継続時間長記憶部から読み出した音韻継続時間長に対応した音素片を選択的に用いて音声合成を行うことを特徴とする音声合成装置。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】 本発明は、定型文章と非定型文章の混在するテキストを音声に変換するような用途に利用し得る音声合成装置に関するものである。

【0002】

【従来の技術】 特開平5-27789号公報には、「音声合成装置」が記載されている。これは、規則合成、分析合成、録音再生など異なる方式の合成音声を作成して用いる場合に、接続部分でオーバーラップ処理を行うものである。特開平8-63187号公報には、定型部分

に対して自然音声から抽出した基本周波数と音韻時間長を用いる「音声合成装置」が開示されている。

【0003】

【発明が解決しようとする課題】 従来、音声合成装置の合成方式には録音編集方式と規則合成方式がある。前者はアナウンサーなどがフレーズ毎に音声を登録しておき、これを適宜選択結合してメッセージを作成するもので、肉声に近い良好な音声を得られる可能性がある反面、データ量が多い、登録外のフレーズには対応できない、新たにフレーズを追加するために同一話者の確保が必要であるなどの問題がある。他方、後者は音素や音節などの細かい単位で音声データを蓄積して任意語彙の合成を可能とするものであるが、音質的に録音編集方式に劣り、特に基本周波数、音韻継続時間長、振幅などの韻律パターンを規則によって付与するためどうしても機械的で不自然なものになる。

【0004】 このため任意語彙の出力が不要な定型的なメッセージには音質の良い録音編集方式が用いられ、テキストからの音声変換が必要な場面では規則合成が用いられる。しかし、カーナビゲーションの音声案内で定型的メッセージの中に地名が埋め込まれるなど、定型文の中に一部任意語彙が埋め込まれるようなアプリケーションも多く存在する。このような場合、ごく一部の任意語彙のために全体の音質を落として規則合成を採用するか、任意語彙の出力を諦めて録音編集方式を用いるか、あるいは定型部分を録音編集で行い任意語彙部分のみ規則合成で行うという混在方式を採るかの選択をせざるを得ない。

【0005】 録音編集と規則合成を混在させる場合の問題点は、2つの方式で出力音声の音質がまったく異なるため、聞いていて違和感があるばかりでなく非常に聞き取り難いものとなる点である。前記特開平5-27789号公報に記載された「音声合成装置」ではこの問題に対し、異なる方式間の出力音声をオーバーラップさせて接続することで対処している。しかし、このようにしても定型部分と任意語部分で話者が変わってしまうことは避けられず基本的な問題は解決していない。また、オーバーラップ部分では2人の話者が同時に話しているようになるため聞き取り難くなる可能性がある。

【0006】 これに対し、前記特開平8-63187号公報に記載された「音声合成装置」では、定型文にも規則合成的に音素あるいは音節等をつないで音韻パラメータを生成し、これに自然音声から抽出した基本周波数及び音韻継続時間長を付与することにより、任意語部分との話者連続性を保持しつつ自然性を向上している。しかし、様々な韻律パラメータや音韻パラメータは相互に関連があり、全体としてバランスを取るよう構築されている規則群の一部（基本周波数と音韻継続時間長）だけを全く異なる話者特性、発声様式の音声から移植することは思わぬ不整合を生んで全体の自然性を損なう可能性

がある。例えば、文末にかけて基本周波数は大きく下がることがあるが、この時は振幅も十分小さくしないと不自然に低い声が目立つことになる。また、このような場合、本来口の開きも小さくなり音声スペクトル自体の変化があるはずであり、あまりに明瞭な音素片データを用いることも違和感を生む。さらに音素片データには対応可能な基本周波数の範囲が存在し、自然音声の基本周波数パターンはこれよりもダイナミックレンジが広いのが普通であるため、無理な基本周波数付与により明瞭性の低下を招く恐れがある。そこで、本発明はこのような問題点を解決し、定型的フレーズの合成音声の自然性を向上することを目的とする。

【0007】

【課題を解決するための手段】請求項1の発明は、人間が発声したフレーズから抽出した音韻継続時間長系列、基本周波数系列、振幅あるいはパワーの系列をそれぞれ記憶した、音韻継続時間記憶部、基本周波数記憶部、振幅記憶部、及び、音素あるいは音素連鎖を音韻単位とし音韻情報を音素片として記憶する音素片記憶部を具備し、入力文字列にしたがって音素片記憶部から読み出し並べた音素片系列を、韻律パターン記憶部から読み出した音韻継続時間長系列にしたがって伸縮して接続し、韻律パターン記憶部から読み出した基本周波数系列にしたがって基本周波数付与を行い、韻律パターン記憶部から読み出した振幅又はパワー系列にしたがって振幅付与を行って音声合成する音声合成装置である。

【0008】請求項2の発明は、請求項1に記載された音声合成装置において、音素片記憶部に記憶した音素片セットによって無理なく合成できる基本周波数の範囲を記憶する基本周波数範囲記憶部を具備し、基本周波数記憶部から読み出した基本周波数系列に対して、基本周波数範囲記憶部から読み出した基本周波数範囲内に収まるように正規化を施して音声合成する音声合成装置である。

【0009】請求項3の発明は、請求項1に記載された音声合成装置において、前記音素片記憶部は一つの音韻単位に対して適応すべき基本周波数範囲毎に複数の音素片を記憶しており、前記基本周波数記憶部から読み出した基本周波数に対応した音素片を選択的に用いて音声合成を行う音声合成装置である。

【0010】請求項4の発明は、請求項1に記載された音声合成装置において、前記音素片記憶部は一つの音韻単位に対して適応すべき振幅範囲毎に複数の音素片を記憶しており、基本周波数記憶部から読み出した振幅あるいはパワーに対応した音素片を選択的に用いて音声合成を行う音声合成装置である。

【0011】請求項5の発明は、請求項1に記載された音声合成装置において、前記音素片記憶部は一つの音韻単位に対して適応すべき音韻継続時間長範囲毎に複数の音素片を記憶しており、音韻継続時間長記憶部から読み

出した音韻継続時間長に対応した音素片を選択的に用いて音声合成を行う音声合成装置である。

【0012】

【発明の実施の形態】本発明の音声合成装置の一実施例について説明する。図1は、この実施例における構成を示す。図1中、M1は音素あるいは音素連鎖を音韻単位とし音韻情報を音素片として記憶する音素片記憶部であって、一つの音韻単位に対して適応すべき基本周波数範囲毎に複数の音素片を記憶しており、基本周波数範囲記憶部M5から読み出した基本周波数に対応した音素片を選択的に用いる。また、M2、M3、M4はそれぞれ人間が発声したフレーズから抽出した音韻継続時間長系列、基本周波数系列、振幅あるいはパワーの系列をそれぞれ記憶した、音韻継続時間長記憶部、基本周波数記憶部、振幅記憶部を示している。M5は基本周波数範囲記憶部であって、音素片記憶部M1に記憶された音素片セットによって無理なく合成できる基本周波数の範囲を記憶している。

【0013】各部の動作について以下に説明する。韻律パターン選択部1は入力される韻律パターン1D（韻律パターンを識別する識別子：例えば、番号等によりそれに対応するパターンを識別するもの）から音韻継続時間長、基本周波数、振幅の各パターンを選択する。音素片選択部2は入力文字列から音素片ラベルを得、また韻律パターン選択部1で選択された音韻継続時間長、基本周波数、振幅の各韻律パターンの範囲を参考にして、これらの情報を元に音素片記憶部M1から必要な音素片を検索する。

【0014】図2は、音素片記憶部M1のデータ構造の一例を示したものである。同一音素ラベル、例えば、“ア”に対して異なる適用可能韻律パラメータ範囲のデータを複数記憶している。ここに示すように韻律パラメータ範囲は、例えば、時間長範囲の長短、増幅範囲の大小などカテゴライズされたものでも、基本周波数範囲のように下限値と上限値を示すものでもよく、またそれらが混在していても構わない。表中のデータ欄には実際には音素片データ、波形データ及びスペクトルパラメータが格納される。この中からラベルの一致するもので韻律パラメータ範囲が最も近いものを選択する。

【0015】音素片伸縮接続部3は、入力文字列にしたがって音素片選択部2により選択された音素片系列を、韻律パターン選択部1で選択された音韻継続時間長の範囲を参考にして、音韻継続時間長記憶部M2から選択された音韻継続時間長パターンに従って伸縮してそれぞれの音素片を接続する。基本周波数範囲記憶部M5には、音素片記憶部M1に記憶された音素片データセット全体によってカバーされる基本周波数範囲が記憶されており、基本周波数パターン正規化部4は、基本周波数範囲記憶部M5から選択された基本周波数パターンがこの範囲を逸脱している場合に、選択された基本周波数パター

ンをこの範囲に合わせて正規化する。基本周波数付与部5は、音素片伸縮接続部3で接続された音素片系列パターンに対して、正規化された基本周波数パターンを付与する。

【0016】振幅付与部6は接続されかつ基本周波数が付与された音素片系列パターンに対し、韻律パターン選択部1で選択された振幅パターンの範囲を参考にして、振幅記憶部M4から選択された振幅パターンを付与して合成音声を作成する。なお、音素片の伸縮及び接続、基本周波数及び振幅の付与に関しては規則音声合成の一般

【0017】

【発明の効果】請求項1に対応する効果：基本的な韻律パラメータである基本周波数、音韻継続時間長、振幅の3つを同じ親善音声フレーズから抽出したものを使用することによって韻律パラメータ間の不整合を抑え、合成音声の自然性を向上することができる。

【0018】請求項2に対応する効果：音素片データベースが対応可能な範囲に基本周波数を正規化することによって、無理な基本周波数付与を防ぎ合成音声の明瞭性の低下を防ぐことができる。

【0019】請求項3に対応する効果：付与すべき基本周波数に対応する音素片データを選択的に用いることにより、ダイナミックレンジの広い自然音声の基本周波数を付与することが可能となり、明瞭性を落とすことなく自然性を向上することが出来る。

【0020】請求項4に対応する効果：同じ音素片デー

タを使い振幅だけを変化させると、スピーカーのボリュームを操作したような機械的な変化となり自然性を損なうが、付与すべき振幅に対応する音素片データを選択的に用いることにより、声の大小によって音韻特性に変化を付けることが出来るため、自然音声のダイナミックレンジの広い振幅変化が付与可能となり、より人間の音声に近い合成音声を得られる。

【0021】請求項5に対応する効果：設定すべき音韻継続時間長に対応する音素片データを選択的に用いることにより、音素片の無理な切りつめによる子音の特徴的部分の欠落や、短い母音定常部分の繰り返しによる機械的な音質を避けることが出来、明瞭性を損なうことなく、自然音声のダイナミックレンジの広いテンポ変化が付与可能となり、より人間の音声に近い合成音声を得られる。

【図面の簡単な説明】

【図1】 本発明による音声合成装置の構成を表すブロック図である。

【図2】 図1に示す音素片記憶部のデータ構造を示したものである。

【符号の説明】

1…韻律パターン選択部、2…音素片選択部、3…音素片伸縮接続部、4…基本周波数パターン正規化部、5…基本周波数付与部、6…振幅付与部、M1…音素片記憶部、M2…音韻継続時間長記憶部、M3…基本周波数記憶部、M4…振幅記憶部、M5…基本周波数範囲記憶部。

【図2】

ラベル	時間長範囲	基本周波数範囲	振幅範囲	データ
ア	長	50Hz~150Hz	大	ア1
ア	短	100Hz~150Hz	大	ア2
ア	標準	50Hz~100Hz	大	ア3
ア	短	80Hz~150Hz	小	ア4
ア	長	30Hz~100Hz	小	ア5
ア	短	30Hz~80Hz	小	ア6
イ	長	50Hz~120Hz	大	イ1
...

【図1】

